



TITLE:

# COMPUTING DOMINANT POLES OF TRANSFER FUNCTIONS (High Performance Algorithms for Computational Science and Their Applications)

AUTHOR(S):

SLEIJPEN, GERARD L.G.; ROMMES, JOOST

---

CITATION:

SLEIJPEN, GERARD L.G. ...[et al]. COMPUTING DOMINANT POLES OF TRANSFER FUNCTIONS (High Performance Algorithms for Computational Science and Their Applications). 数理解析研究所講究録 2008, 1614: 111-123

ISSUE DATE:

2008-10

URL:

<http://hdl.handle.net/2433/140109>

RIGHT:

## COMPUTING DOMINANT POLES OF TRANSFER FUNCTIONS

GERARD L.G. SLEIJPEN\* AND JOOST ROMMES†

**Abstract.** The transfer function describes the response of a dynamical system to periodic inputs. Dominant poles are specific eigenvalues of the state space matrix that corresponds to the dynamical system. The number of these type of eigenvalues is often small as compared to the number of state variables (the dimension of the state space matrix). The dominant behaviour of the dynamical system can be captured by projecting the state-space onto the subspace spanned by the eigenvectors corresponding to dominant poles.

In this paper, we discuss numerical methods for computing these poles and corresponding eigenvectors.

**Keywords:** Transfer function, dominant pole, dynamical system, eigenvectors, Rayleigh quotient iteration, Jacobi-Davidson.

**AMS(MOS) subject classification:** 65F15, 65N25.

**1. Introduction.** We are interested in methods for evaluating the *transfer function*

$$(1.1) \quad H(\omega) \equiv \mathbf{c}^*(2\pi i\omega \mathbf{E} - \mathbf{A})^{-1} \mathbf{b} \quad (\omega \in \mathbb{R}),$$

where  $\mathbf{A}$  and  $\mathbf{E}$  are given real  $N \times N$  matrices and  $\mathbf{b}$  and  $\mathbf{c}$  are given real  $N$ -vectors. The function  $H$  plays an important role in the analysis of the dynamical behavior of constructions (linear systems) as buildings, airplanes, electrical circuits (chips), and of phenomena as tidal movements in oceans and bays, etc. (see §2.1).

In practice, the problem of computing  $H$  is a computational challenge. The dimension is often high ( $N$  is large), the matrices  $\mathbf{A}$  and  $\mathbf{E}$  are often sparse, but they can be unstructured. Although function values  $H(\omega)$  can be obtained by solving linear systems, this approach is usually not fruitful: the function values have to be computed for a wide range of  $\omega$  (say, between 0 and  $10^4$ ), and a preconditioner that is effective for  $2\pi i\omega \mathbf{E} - \mathbf{A}$  for one value of  $\omega$  will not work for another value.

In this paper, we will approximate the *system*  $(\mathbf{A}, \mathbf{E}, \mathbf{b}, \mathbf{c})$  by a (much) smaller one by projection onto spaces of eigenvectors associated to specific eigenvalues of the pencil  $(\mathbf{A}, \mathbf{E})$ , to the so-called dominant poles. The number of dominant poles is usually much smaller than  $N$  (see §2.5) and, therefore, this *modal approach* can be very effective.

We discuss numerical methods as the dominant pole algorithm (DPA [15], see §3.1) and Rayleigh quotient iteration (RQI [19, 20], see §3.3) for computing these poles and corresponding eigenvectors. Although DPA, is based on an old variant of inverse iteration [18], we argue that it has more attractive convergence properties for finding dominant poles than the celebrated RQI (see §3.5). RQI as well as DPA iterates on single vectors. To accelerate convergence, extensions of these methods that built search subspace, lead to variants of the Jacobi-Davidson [28] method in §5.

The fact that a Krylov subspace generated by a matrix  $\tilde{\mathbf{A}}$  and a vector  $\tilde{\mathbf{b}}$  equals the Krylov subspace generated by  $\mathbf{I} - 2\pi i(\omega_0 - \omega)\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{b}}$  suggests that (unpreconditioned) Krylov solvers that work well for one value of  $\omega$  might also be effective for other values of  $\omega$  (see §6 for more details). Moment-matching techniques (see, e.g., [11, 1]) are based on this observation. These techniques tend to produce approximations that are highly accurate in the neighborhood of  $\omega_0$ . Variants as the rational Krylov sequence (RKS, [25]) method are efficient and allow good approximations around other values of  $\omega$  as well. These approximations tend to be accurate on the “smooth” part of  $H$ , whereas the modal approach is highly accurately around those values of  $\omega$ , where  $H$  has ‘peaks’. In §6, we will see that combining both approaches improves accuracy and efficiency.

This paper summarizes the first half of [22].

In this paper, we use the Euclidean norm denoted by  $\|\cdot\|$ .

## 2. Transfer functions and dominant poles.

\*Mathematical Institute, Utrecht University, P.O. Box 80010, 3508 TA Utrecht, the Netherlands, e-mail: sleijpen@math.uu.nl

†NXP Semiconductors Corporate I&T / Design Technology & Flows, PostBox WY4-01, NL-5656 AE Eindhoven, The Netherlands, e-mail: joost.rommes@nxp.com

**2.1. Linear, time invariant, dynamical systems.** Linear, time invariant, dynamical system can often be modelled as (see, e.g., [5])

$$(2.1) \quad \begin{cases} \mathbf{E} \dot{\mathbf{x}}(t) &= \mathbf{A} \mathbf{x}(t) + \mathbf{b} u(t) \\ y(t) &= \mathbf{c}^* \mathbf{x}(t) + d u(t). \end{cases}$$

Here,  $\mathbf{A}$  and  $\mathbf{E}$  are  $N \times N$  matrices,  $\mathbf{b}$ ,  $\mathbf{c}$  and  $d$  are matrices of size  $N \times m$ ,  $N \times p$ , and  $p \times m$ , respectively,  $u$ ,  $\mathbf{x}$  and  $y$  are functions of  $t \in \mathbb{R}$ , with values in  $\mathbb{R}^m$ ,  $\mathbb{R}^N$ ,  $\mathbb{R}^p$ , respectively. The matrix  $\mathbf{E}$  may be non-singular, but we will assume that  $(\mathbf{A}, \mathbf{E})$  is a regular pencil, that is,  $s \mapsto \det(s\mathbf{E} - \mathbf{A})$  is non trivial on  $\mathbb{C}$ . The *system*  $(\mathbf{A}, \mathbf{E}, \mathbf{b}, \mathbf{c}, d)$  is given. The *control function*  $u$  determines the *state vector*  $\mathbf{x}$  and, the function of interest, the *output*  $y$  of the system.  $\mathbf{A}$  is the *state matrix*,  $N$  is the *number of states* or *order of the system*,  $\mathbf{b}$  is the *input map* of the system, and  $\mathbf{c}$  is the *output map*.

In this paper, we restrict our discussion to the case where  $m = p = 1$ , a SISO (single input single output) system. In practice, however, one usually has to deal with MIMO (multiple input multiple output) systems ( $m > 1$ ,  $p > 1$ ) (cf., e.g., [23]). Moreover, the systems are often non-linear. Then, the linear systems arise by linearization (in a Newton process). The number of states is large (in the range of  $10^4 - 10^8$ ), the matrices  $\mathbf{A}$  and  $\mathbf{E}$  are sparse and (often) unstructured.

In this paper, bold face letters refer to high dimensional quantities.

Dynamical systems play an important role in, for instance, the stability analysis of technical constructions as airplanes, buildings, bridges, etc.. The behavior of such constructions can be modelled by a set of partial differential equations using laws from structural mechanics. Discretization of the spatial part leads to high dimensional ordinary differential equations as in the first equation of (2.1). Motions of the constructions can be induced by dynamical forces acting on certain points. These action are modelled by the term  $\mathbf{b} u(t)$ . The second equation of (2.1) models the response at certain (other) points of the construction.

For instance, a building can shake in an earthquake. The earthquake applies a (periodic) force at the foundations of the building. The resulting swing at the toplevel may be of interest. To guarantee that the building survives earthquakes, the swing at, for instance, the top level should be small. To achieve this, the design of the building may have to be adapted, which means that (2.1) has to be solved for a different system. Of course, design and computations have to be performed before the actual construction of the building.

In electrical circuit simulations, the matrices  $\mathbf{A}$  and  $\mathbf{E}$  incorporate the incidence matrix of the directed graph that describes the circuit. The values of the electrical components (as resistors, capacitors, ...) determine the value of matrix coefficients. These matrices are unstructured and the values of (neighboring) matrix coefficients may vary orders of magnitude.

**2.2. Transfer functions.** To analyze system (2.1), apply a control function  $u$  of the form  $u(t) = \exp(st)$  ( $t \in \mathbb{R}$ ) for some  $s \in \mathbb{C}$ , that is, apply a Laplace transform. Then  $\mathbf{x}(t) = (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{b} \exp(st)$  and  $y(t) = [\mathbf{c}^*(s\mathbf{E} - \mathbf{A})^{-1} \mathbf{b} + d] \exp(st)$ .

The function

$$(2.2) \quad H(s) \equiv \mathbf{c}^*(s\mathbf{E} - \mathbf{A})^{-1} \mathbf{b} + d \quad (s \in \mathbb{C})$$

is the *transfer function* of the system (2.1). The transfer function describes to response of the system at the output to a (damped) periodic force at the input. The response to harmonic oscillations, i.e.,  $s = 2\pi i\omega$  with *frequency*  $\omega$  in  $\mathbb{R}$ , is of particular interest.

Since the computational complexity is not affected by  $d$ , we further assume that  $d = 0$ .

**2.3. Model order reduction.** Note that, in principle,  $H$  can be computed by only solving linear systems. Nevertheless, computing the transfer function is often very hard. In practice,  $N$  is large,  $H$  is needed for a wide range of  $\omega$  (in  $\mathbb{R}$ ,  $s = 2\pi i\omega$ ), preconditioners are hard to include (a preconditioner for  $\mathbf{A}$  is not a preconditioner for  $s\mathbf{E} - \mathbf{A}$  for  $s \neq 0$ ), and, specifically in a design stage, (or in case the system is non-linear) solutions are required for a number of (slightly different) systems.

For these reasons, *reduced model* are constructed, that is,  $k$ th order systems  $(\tilde{\mathbf{A}}, \tilde{\mathbf{E}}, \tilde{\mathbf{b}}, \tilde{\mathbf{c}})$  with  $k \ll N$  for which (cf., [2])

- for all  $u$ , the error  $\|y(t) - \tilde{y}(t)\|$  is ‘small’ in some norm (as the 2-norm, Hankel-norm, ...),<sup>1</sup>
- physical and numerical properties (as, stability, passivity, ...) are preserved;
- the reduced system can be efficiently and stably computed;
- the error can be efficiently measured when constructing a reduced system.

In practice, there are other restrictions as well. System (2.1) may have a certain (block) structure. For instance, if the first order differential equation in (2.1) stems from a second order one. Then  $\mathbf{x}$  has a block of first order derivative terms. It may be convenient to preserve this (block) structure (cf., [3]). Designers would like to see reduced models that are realizable. In, for instance, electrical circuit design, they would like to see reduced models that can be build as an electrical circuit. Also reduced models require computational efforts. Therefore, it would be helpful if these smaller models also would ‘fit’ in existing simulation software.

**2.4. Approaches for model order reduction.** Transfer functions can be expressed in terms of eigenvalues and eigenvectors as we will see in §2.5 below. Approaches for constructing reduced models correspond to approaches for computing eigenvalues and eigenvectors. Three main classes of approaches can be distinguished.

1) Methods based on balanced truncation [16] and Hankel norm approximation [9]. They form the analogue of the QR and QZ methods for computing eigenvalues.

2) Padé approximation and moment-matching techniques [10]. They are usually based on Krylov subspace techniques as (shift-and-invert) bi-Lanczos and Arnoldi (see §6.1). There are block versions, two sided versions and versions based on rational Krylov sequence [25].

3) Modal approximations [7], where, as we will see in §2.5, the reduced model is based on dominant poles. Such approaches correspond to Jacobi–Davidson type of methods for computing a few eigenvalues and eigenvectors.

In the second and third type of approach, matrices  $\mathbf{V}_k$  and  $\mathbf{W}_k$  of size  $n \times k$  are computed, and the reduced model arises by projection:

$$(2.3) \quad \tilde{\mathbf{A}} = \mathbf{W}_k^* \mathbf{A} \mathbf{V}_k, \quad \tilde{\mathbf{E}} = \mathbf{W}_k^* \mathbf{E} \mathbf{V}_k, \quad \tilde{\mathbf{b}} = \mathbf{W}_k^* \mathbf{b}, \quad \tilde{\mathbf{c}} = \mathbf{V}_k^* \mathbf{c}.$$

In bi-Lanczos and Arnoldi, the columns of  $\mathbf{V}_k$  form a basis of the Krylov subspace  $\mathcal{K}_k((s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{E}, (s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{b})$  for some  $s_0 \in \mathbb{C}$ . In the modal approaches, the columns of  $\mathbf{V}_k$  and  $\mathbf{W}_k$  span appropriate right and left, respectively, eigenvectors of the pair  $(\mathbf{A}, \mathbf{E})$ .

In this paper, we mainly concentrate on the modal approach.

**2.5. Dominant poles.** Before we explain how dominant poles play an important role in our approach, we first define a pole and introduce the concept of dominance.

A  $\lambda \in \mathbb{C}$  is a *pole* of  $H$  if  $\lim_{\mu \rightarrow \lambda} |H(\mu)| = \infty$ .

Poles are eigenvalues of the pencil  $(\mathbf{A}, \mathbf{E})$ : non-zero  $\mathbf{v}_i$  and  $\mathbf{w}_i$  are right and left, respectively, eigenvectors with eigenvalue  $\lambda_i$  if

$$(2.4) \quad \mathbf{A} \mathbf{v}_i = \lambda_i \mathbf{E} \mathbf{v}_i \quad \text{and} \quad \mathbf{w}_i^* \mathbf{A} = \lambda_i \mathbf{w}_i^* \mathbf{E}.$$

We assume that  $(\mathbf{A}, \mathbf{E})$  is nondefective. In each right eigenspace corresponding to one eigenvalue, we select the one eigenvector that determines the component of  $\mathbf{b}$  in that space. Similarly, for  $\mathbf{c}$  and left eigenvectors. To be more precise, we select the *eigen triples*  $(\mathbf{v}_i, \mathbf{w}_i, \lambda_i)$  such that  $\lambda_i \neq \lambda_j$  ( $i \neq j$ ),

$$\mathbf{b} = \sum_{i=1}^n \beta_i \mathbf{E} \mathbf{v}_i + \beta_\infty \mathbf{A} \mathbf{v}_\infty, \quad \mathbf{c} = \sum_{i=1}^n \gamma_i \mathbf{E}^* \mathbf{w}_i + \gamma_\infty \mathbf{A}^* \mathbf{w}_\infty,$$

with  $\mathbf{v}_\infty$  and  $\mathbf{w}_\infty$  in the kernel of  $\mathbf{E}$  and  $\mathbf{E}^*$ , respectively. In addition, we assume the eigenvectors to be scaled such that  $\mathbf{w}_i^* \mathbf{E} \mathbf{v}_i = 1$  if  $\mathbf{w}_i^* \mathbf{E} \mathbf{v}_i \neq 0$ . Now, note that [13],

$$(2.5) \quad H(s) = \sum_{i=1}^n \frac{R_i}{s - \lambda_i} + R_\infty, \quad \text{where} \quad R_i \equiv (\mathbf{c}^* \mathbf{v}_i)(\mathbf{w}_i^* \mathbf{b}) = \bar{\gamma}_i \beta_i.$$

<sup>1</sup>We apologize for using terms as ‘Hankel norm’, ‘stability’ and ‘passivity’ that have not been defined. These term do not play an essential role in this paper.

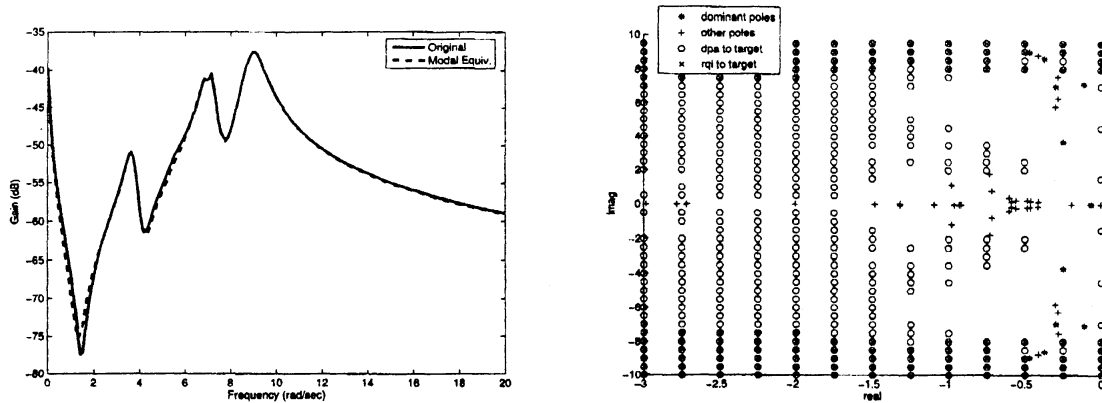


FIG. 2.1: The left figure show that Bode plot of the New England test system [15] (solid line) together with the approximation (dashed line) from the modal approach using the 11 most dominant poles out of 66 eigenvalues ( $N = 66$ ).

The right figure shows part of the complex plane with part of the spectrum (dominant poles are marked with \*), together with a grid of location for the initial shifts  $s_0$  in the complex plane for which DPA (marked with o) and for which RQI (marked with x) converge to the most dominant pole  $\lambda = -0.476 \pm 8.96i$ . Initial shifts  $s_0$  on the displayed grid without marker lead to convergence to less dominant poles.

The  $R_i$  are called *residues*. Note that  $R_i = 0$  if  $\mathbf{w}_i^* \mathbf{E} \mathbf{v}_i = 0$ .

The ‘contribution’ of the pole  $\lambda_i$  to the transfer function is determined by the size  $|R_i|$  of the corresponding residue and the distance  $|\text{Re}(\lambda_i)|$  if the pole to the imaginary axis: the pole  $\lambda_i$  is said to be *dominant* if the scaled size  $\frac{|R_i|}{|\text{Re}(\lambda_i)|}$  of the corresponding residue is large.

Dominant poles determine high peaks in the so-called *Bode plot*, that is, in the graph of  $\omega \rightsquigarrow |H(2\pi i \omega)|$  ( $\omega \in \mathbb{R}$ ). For an example, see the left panel of Fig. 2.1, where this graph is plotted on decibel scale. The value  $|H(2\pi i \omega)|$  is the *gain* of the system at frequency  $\omega$ .

In the *modal approach*, the model is reduced by projecting onto the space spanned by eigenvectors corresponding to the poles that are most dominant,  $\mathbf{b}$  onto right eigenvector and  $\mathbf{c}$  onto corresponding left eigenvectors. Note that, our scaling implies that  $\bar{\mathbf{E}} = \mathbf{I}$  (see (2.3) and use that  $\mathbf{w}_j^* \mathbf{E} \mathbf{v}_i = 0$  if  $\lambda_i \neq \lambda_j$ ).

The success of the modal approach comes from the fact that, in practice, the number of dominant poles is usually much smaller than the number of poles (which is smaller than number  $n$  of relevant eigenvalues, while  $n \leq N$ ).

**3. Algorithms and convergence.** Dominance of poles has been defined (see §2.5) in a relative and a subjective sense (‘large’), relative with respect to the distance from the imaginary axis. For mathematical reasons, however, we will in this section assume dominance in an absolute sense [12]:  $\lambda_i$  is *dominant* if  $|R_i| > |R_j|$  for all  $j \neq i$ .

**3.1. Dominant pole algorithm.** The (generalized two-sided) Rayleigh quotient [19, 20]

$$\rho(\mathbf{x}, \mathbf{y}) \equiv \rho(\mathbf{A}, \mathbf{E}, \mathbf{x}, \mathbf{y}) \equiv \frac{\mathbf{y}^* \mathbf{A} \mathbf{x}}{\mathbf{y}^* \mathbf{E} \mathbf{x}} \quad (\text{provided } \mathbf{y}^* \mathbf{E} \mathbf{x} \neq 0),$$

will play an important role in the algorithms. Note that  $\mathbf{y}^* \mathbf{E} \mathbf{x}$  can be 0 even if  $\mathbf{E}$  is non-singular. If  $\mathbf{y} = \mathbf{x}$ , we put  $\rho(\mathbf{x})$ :  $\rho(\mathbf{x}) \equiv \rho(\mathbf{x}, \mathbf{x})$ .

Since poles are zeros of the function  $s \rightsquigarrow 1/H(s)$ , Newton’s process can be applied:

$$s_{k+1} = s_k - \frac{\mathbf{c}^*(s_k \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}}{\mathbf{c}^*(s_k \mathbf{E} - \mathbf{A})^{-1} \mathbf{E} (s_k \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}}.$$

With  $\mathbf{x}_k \equiv (s_k \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}$  and  $\mathbf{y}^* \equiv \mathbf{c}^*(s_k \mathbf{E} - \mathbf{A})^{-1}$ ,  $s_{k+1}$  can be expressed as a Rayleigh quotient:  $s_{k+1} = \rho(\mathbf{x}_k, \mathbf{y}_k)$ . This leads to the algorithm (cf. [15]) represented in the left panel of Alg. 3.1, where the indices have been suppressed to indicate that new values can replace old ones (indicated by the same letter). The algorithm is represented in its simplest form: stopping

<p>Select <math>s_0 \in \mathbb{C}</math> and <math>tol &gt; 0</math>.  Set <math>\nu = 1, s = s_0</math></p> <p>While <math>\nu &gt; tol</math> repeat</p> <p style="padding-left: 40px;">Solve <math>(s\mathbf{E} - \mathbf{A})\mathbf{x} = \mathbf{b}</math> for <math>\mathbf{x}</math>  Solve <math>(s\mathbf{E} - \mathbf{A})^*\mathbf{y} = \mathbf{c}</math> for <math>\mathbf{y}</math>  <math>s = \frac{\mathbf{y}^*\mathbf{A}\mathbf{x}}{\mathbf{y}^*\mathbf{E}\mathbf{x}}</math>  <math>\nu = \ \mathbf{A}\mathbf{x} - s\mathbf{E}\mathbf{x}\ </math></p> <p>end while</p>	<p>Select <math>s_0 \in \mathbb{C}</math> and <math>tol &gt; 0</math>.  Set <math>\nu = 1, s = s_0</math>  <math>\mathbf{x} = (s\mathbf{E} - \mathbf{A})^{-1}\mathbf{b}, \mathbf{y} = (s\mathbf{E} - \mathbf{A})^{-*}\mathbf{c}</math></p> <p>While <math>\nu &gt; tol</math> repeat</p> <p style="padding-left: 40px;"><math>\tilde{\mathbf{b}} = \mathbf{E}\mathbf{x}/\ \mathbf{x}\ , \tilde{\mathbf{c}} = \mathbf{E}^*\mathbf{y}/\ \mathbf{y}\ </math>  Solve <math>(s\mathbf{E} - \mathbf{A})\mathbf{x} = \tilde{\mathbf{b}}</math> for <math>\mathbf{x}</math>  Solve <math>(s\mathbf{E} - \mathbf{A})^*\mathbf{y} = \tilde{\mathbf{c}}</math> for <math>\mathbf{y}</math>  <math>s = \frac{\mathbf{y}^*\mathbf{A}\mathbf{x}}{\mathbf{y}^*\mathbf{E}\mathbf{x}}</math>  <math>\nu = \ \mathbf{A}\mathbf{x} - s\mathbf{E}\mathbf{x}\ </math></p> <p>end while</p>
---	---

**Alg. 3.1:** The left panel represents the dominant pole algorithm, while the right panel represents Rayleigh quotient iteration. An approximate eigentriple is formed by  $(\mathbf{x}, \mathbf{y}, s)$ .  $\nu$  is the residual of the approximate right eigenvector.

criteria that are more sophisticated (and more appropriate) than ' $\nu_k \equiv \|\mathbf{A}\mathbf{x}_k - s_{k+1}\mathbf{E}\mathbf{x}_k\| < tol$ ' can be exploited as well.

We refer to this algorithm as the *dominant pole algorithm* (DPA), since, as we will argue below, this algorithm tends to find dominant poles first.

The algorithm is applicable in practice if it is feasible to form the LU-decomposition of  $s\mathbf{E} - \mathbf{A}$ . We will come back to this restriction in §5.5.

**3.2. Convergence.** DPA represents a Newton process. Therefore, the sequence  $(s_k)$  will converge quadratically to an eigenvalue  $\lambda$  provided that the initial value  $s_0$  is sufficiently close to  $\lambda$ . The  $\mathbf{x}_k$  and  $\mathbf{y}_k$  form approximate eigenvectors and they converge quadratically as well:

**THEOREM 3.1** ([24], Th. 4.2). *If  $(\mathbf{v}, \mathbf{w}, \lambda)$  is an eigentriple and DPA has been applied to produce  $(\mathbf{x}_k, \mathbf{y}_k, s_{k+1})$ , then*

$$\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{v} \Leftrightarrow \lim_{k \rightarrow \infty} \mathbf{y}_k = \mathbf{w} \Leftrightarrow \lim_{k \rightarrow \infty} s_k = \lambda.$$

*Convergence implies quadratic convergence: for some  $\kappa > 0$ , we have*

$$\|\mathbf{v} - \mathbf{x}_{k+1}\| \leq \kappa \|\mathbf{v} - \mathbf{x}_k\| \|\mathbf{w} - \mathbf{y}_k\| \quad \text{and} \quad \|\mathbf{w} - \mathbf{y}_{k+1}\| \leq \kappa \|\mathbf{v} - \mathbf{x}_k\| \|\mathbf{w} - \mathbf{y}_k\|$$

**3.3. Rayleigh quotient iteration.** DPA uses the best available eigenvector approximations to update the approximate eigenvalue:  $s_{k+1} = \rho(\mathbf{x}_k, \mathbf{y}_k)$ . The approximate eigenvectors can also be exploited in the computation of the new approximate eigenvectors:  $\mathbf{x}_k = (s_k\mathbf{E} - \mathbf{A})^{-1}\mathbf{E}\mathbf{x}_{k-1}$ , and  $\mathbf{y}_k^* = \mathbf{y}_{k-1}^*\mathbf{E}(s_k\mathbf{E} - \mathbf{A})^{-1}$ . This leads to the celebrated (two-sided) Rayleigh quotient iteration (RQI [19, 20], see the right panel of Alg. 3.1) for which faster convergence is to be expected: of RQI it is known that convergence implies cubic convergence [20, p.689].

**3.4. Initiation.** Both methods, DPA as well as RQI, require initiation. The selection  $s_0 = \rho(\mathbf{b}, \mathbf{c})$  in DPA may seem reasonable. It corresponds to a choice of  $\mathbf{x}_0 = \mathbf{b}$  and  $\mathbf{y}_0 = \mathbf{c}$  in RQI. In the symmetric case, where  $\mathbf{A}^* = \mathbf{A}$ ,  $\mathbf{E} = \mathbf{I}$  and  $\mathbf{c} = \mathbf{b}$ , this choice works well. Unfortunately, in the general case,  $\mathbf{c}^*\mathbf{E}\mathbf{b}$  can be zero (or very small), and, in practice, we observed that it often is. Then  $\rho(\mathbf{b}, \mathbf{c})$  is not defined (or very large). Therefore, we select an  $s_0$  and we proceed as indicated in the algorithms in Alg. 3.1.

**3.5. Convergence regions.** Both methods, DPA and RQI, converge fast (if they converge). However, we are not just interested in (quickly) finding eigentriples. We want the *dominant* ones, that is, the dominant poles with associated eigenvectors. Since, local convergence is guaranteed, both algorithms are able to detect these, but detection depends on how close the initial  $s_0$  is to the wanted dominant poles. The question is what is 'close'? The *convergence region* of a pole  $\lambda$  and a method is the collection of  $s_0$  in  $\mathbb{C}$  for which the sequence  $(s_k)$

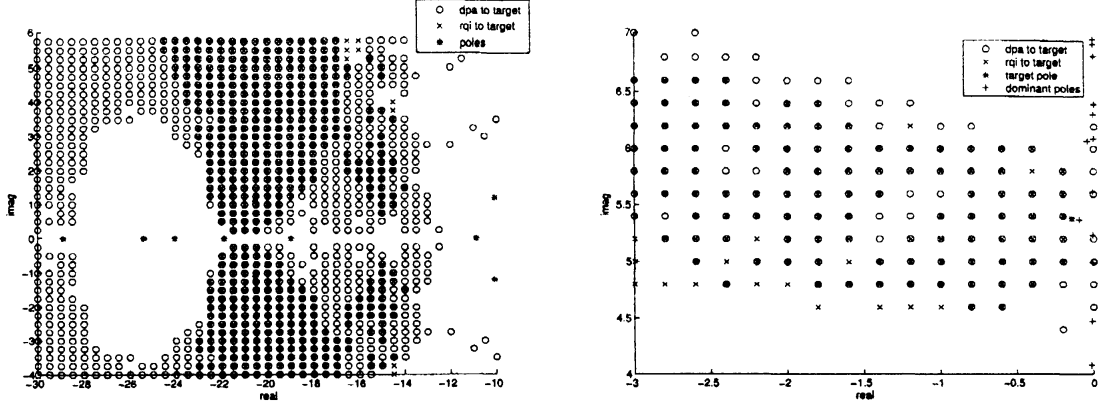


FIG. 3.1: The figures show convergence regions for DPA and RQI for two examples (Ex. 3.2). The most dominant pole in the displayed part of the complex plain is at  $\lambda \approx -20.5 \pm 1.1i$  in the left panel and at  $\lambda \approx -0.14 \pm 5.4i$  in the right panel. For a description of the symbols used in the figures, see Fig. 2.1.

produced by the method converges towards  $\lambda$ . Experiments (as in the right panel of Fig. 2.1 and in the examples below) indicate that the region of convergence of the dominant poles is much larger for DPA than for RQI. This advantage compensates for a slightly slower convergence (recall that both methods locally converge fast and it is not clear in advance whether in practice there is any advantage associated to cubic convergence over quadratic convergence. We observed that RQI typically needed only 10%-20% less iteration steps than DPA).

Fig. 3.1 displays the results for following two examples (see [24]).

EXAMPLE 3.2. a. The left panel displays the convergence regions for the “Brazilian Interconnect Power System” [21]. The order of this test model is  $n = 13,251$ . The matrix  $\mathbf{E}$  is singular,  $\mathbf{b}$  and  $\mathbf{c}$  only have one nonzero entry and  $\mathbf{c}^* \mathbf{E} \mathbf{b} = 0$ . The most dominant poles appear in complex conjugate pairs. From the figure, we learn that for many initial shifts DPA converges to the most dominant pole, where RQI does not, while for a small number of initial shifts, RQI converges to the most dominant pole where DPA does not.

b. The second test model, with results displayed in the right panel, is the PEEC system [6], a well-known benchmark system for model order reduction. One of the difficulties with this system of order  $n = 480$  is that it has many equally dominant poles that lie close to each other in a relatively small part  $[-1, 0] \times [-10, 10]i$  in the complex plane. Although the difference is less pronounced than the previous example, DPA converges to the most dominant pole in more cases than RQI.

The examples show that the convergence region of the dominant pole for DPA is much larger than for RQI. For an heuristic explanation, recall that dominance of a pole  $\lambda_i$  is determined by the size of the associated residue  $R_i = (\mathbf{c}^* \mathbf{v}_i)(\mathbf{w}_i^* \mathbf{b})$ , which depends on both  $\mathbf{b}$  and  $\mathbf{c}$ . DPA uses information from  $\mathbf{b}$  and  $\mathbf{c}$  in *each* step (recall that  $\mathbf{x}_k = (s_k \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}, \dots$ ), whereas  $\mathbf{b}$  and  $\mathbf{c}$  enter the RQI process in the initiation phase only. RQI converges fast, but it tends to converge towards the eigenvalue closest to  $s_0$  rather than to the one for which  $|\mathbf{c}^* \mathbf{v}_i|$  and  $|\mathbf{w}_i^* \mathbf{b}|$  are large.

In the next subsection, we will see that, for the symmetric case, there are mathematical arguments that underline this explanation.

**3.6. Convergence in the symmetric case.** In this subsection, we assume that  $\mathbf{A}^* = \mathbf{A}$ ,  $\mathbf{E} = \mathbf{I}$ , and  $\mathbf{c} = \mathbf{b}$ , and will concentrate on the eigenpair  $(\mathbf{v}, \lambda)$ :  $\mathbf{A} \mathbf{v} = \lambda \mathbf{v}$ .

The *gap*  $\gamma \equiv \min\{|\lambda_i - \lambda| \mid \lambda_i \neq \lambda\}$  between  $\lambda$  and the other eigenvalues  $\lambda_i$  of  $\mathbf{A}$  will play a role and  $\zeta \equiv |\tan \angle(\mathbf{v}, \mathbf{b})|$ , which expresses the angle between the eigenvector  $\mathbf{v}$  and  $\mathbf{b}$ . The initial  $s_0$  is selected in  $\mathbb{R}$ .

If  $(\mathbf{x}_k)$  is a sequence of approximate eigenvector, then, for  $k \geq 0$ , we put

$$s_{k+1} \equiv \frac{\mathbf{x}_k^* \mathbf{A} \mathbf{x}_k}{\mathbf{x}_k^* \mathbf{x}_k}, \quad \alpha_k \equiv \frac{|s_k - \lambda|}{\gamma - |s_k - \lambda|}, \quad \zeta_{k+1} \equiv \tan \angle(\mathbf{v}, \mathbf{x}_k).$$

THEOREM 3.3 ([24], Th.4.3-5). a. If  $(\mathbf{A} - s_k \mathbf{I})\mathbf{x}_k = \mathbf{b}$  (DPA), then

$$\alpha_0 < \alpha_{\text{DPA}} \equiv \frac{1}{\zeta^2} \quad \text{implies that} \quad s_k \rightarrow \lambda \quad \text{and} \quad \alpha_{k+1} \zeta^2 \leq (\alpha_k \zeta^2)^2 < 1.$$

b. If  $\mathbf{x}_{-1} = \mathbf{b}$  and  $(\mathbf{A} - s_k \mathbf{I})\mathbf{x}_k = \mathbf{x}_{k-1}$  (RQI), then

$$\alpha_0 < \alpha_{\text{RQI}} \equiv \frac{1}{\zeta} \quad \text{implies that} \quad s_k \rightarrow \lambda \quad \text{and} \\ \alpha_{k+1} \leq (\alpha_k \zeta_k)^2, \quad \zeta_{k+1} \leq \alpha_k \zeta_k, \quad \alpha_{k+1} \zeta_{k+1} \leq (\alpha_k \zeta_k)^3 < 1.$$

The last estimate in Th.3.3.b expresses cubic convergence. Note that this estimate is the product of the two preceding ones. The estimate in Th.3.3.a expresses quadratic convergence.

The theorem is of interest for two reasons. The use of a scaled error  $\alpha_k$  in the eigenvalue approximation leads to the elegant estimates of Th.3.3, and the estimates are sharp (for any symmetric matrix  $\mathbf{A}$ , there is a vector  $\mathbf{b}$  such that for any  $s_0$  between  $\lambda$  and  $\lambda_{i_0}$ ,  $(s_k)$  converges to  $\lambda$  only if the condition in the theorem is satisfied. Here,  $i_0$  is such that  $\gamma = |\lambda - \lambda_{i_0}|$ .) and sharper than previous results in literature: for instance, for DPA, see [18], and for RQI, see [4].

Note that  $\zeta < 1$  implies that  $\lambda$  is dominant, because the residue associated with  $\lambda$  is  $|\mathbf{v}^* \mathbf{b}|^2 > \frac{1}{2} \|\mathbf{b}\|^2$ , whence the residues for the other eigenvalues are  $< \frac{1}{2} \|\mathbf{b}\|^2$ . Moreover, if  $\zeta < 1$ , then  $1/\zeta < 1/\zeta^2$  and the theorem tells us that the convergence region for the dominant eigenvalue  $\lambda$  is larger for DPA than for RQI.

**4. Deflation.** For the modal approach, eigentriples with poles that are most dominant are required. Suppose one (dominant) eigentriple  $(\mathbf{v}, \mathbf{w}, \lambda)$  has been detected. Then, one can select a new initial shift  $s_0$ , hoping that this leads to another (dominant) eigentriple. But it is not unlikely that the same eigentriple will be detected. A more reliable strategy for avoiding this situation is deflation, where the detected eigentriple is removed from the process of finding dominant eigentriples.

DPA allows efficient deflation as follows (see, [21, Th.3.1])

THEOREM 4.1. With  $\tilde{\mathbf{b}} \equiv (\mathbf{I} - \mathbf{E} \mathbf{v} \mathbf{w}^*) \mathbf{b}$  and  $\tilde{\mathbf{c}}^* \equiv \mathbf{c}^* (\mathbf{I} - \mathbf{v} \mathbf{w}^* \mathbf{E})$ , we have that

$$\tilde{H}(s) \equiv \tilde{\mathbf{c}}^* (s \mathbf{E} - \mathbf{A})^{-1} \tilde{\mathbf{b}}$$

has the same poles and residues as  $H(s) = \mathbf{c}^* (s \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}$ , except for the residue associated with  $\lambda$  which is transformed to 0.

Since the ‘new’ residue is zero,  $\lambda$  is not dominant for the system  $(\mathbf{A}, \mathbf{E}, \tilde{\mathbf{b}}, \tilde{\mathbf{c}})$ , while for all other poles the residues are unchanged. In particular, DPA can be applied to this ‘deflated’ system to detect the eigentriple that is next in dominance. Of course, this deflation strategy can be repeated to find the third dominant eigentriple, etc..

This deflation is cheap: it has to be applied to two vectors only, to  $\mathbf{b}$  and  $\mathbf{c}$  for finding the second eigentriple, to  $\tilde{\mathbf{b}}$  and  $\tilde{\mathbf{c}}$  for finding the third triple, etc..

Theoretically, the same deflation strategy can be applied in combination with RQI (and to other methods, as Jacobi-Davidson, for computing eigentriples). Unfortunately, in practice, due to rounding errors, this efficient strategy will not prevent RQI from finding the same eigentriple. To see this, recall that the inverse power method

$$\mathbf{x}_k = \tilde{\mathbf{x}}_k / \|\tilde{\mathbf{x}}_k\|, \quad \tilde{\mathbf{x}}_{k+1} = (s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{E} \mathbf{x}_k$$

is likely to converge to the eigentriple  $(\mathbf{v}_j, \mathbf{w}_j, \lambda_j)$  with  $\lambda_j$  such that  $1/(s_0 - \lambda_j)$  is the absolute largest eigenvalue of  $(s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{E}$ . When  $\mathbf{v}_j$  is deflated from an initial guess  $\mathbf{x}_0$  (by multiplication by  $\mathbf{I} - \mathbf{v}_j \mathbf{w}_j^* \mathbf{E}$ ),<sup>2</sup> then the next iterate  $\mathbf{x}_1$  will have a component in the direction of  $\mathbf{v}_j$  due to rounding errors. Though the component will be small, it will be amplified in the next iterations and will eventually lead to the same eigentriple  $(\mathbf{v}_j, \mathbf{w}_j, \lambda_j)$ : the inverse power method blows up small unwanted components. To prevent this from happening, the deflation

<sup>2</sup>Note that, for  $\mathbf{E} \mathbf{x}_0$ , this is equivalent to multiplying  $\mathbf{E} \mathbf{x}_0$  by  $\mathbf{I} - \mathbf{E} \mathbf{v}_j \mathbf{w}_j^*$ , which is the projection that deflates  $\mathbf{b}$  in Th. 4.1.



(multiplication by  $\mathbf{I} - \mathbf{v}_j \mathbf{w}_j^* \mathbf{E}$ ) has to be performed to each of the iterates  $\mathbf{x}_k$ , which makes the deflation relatively expensive. A similar remark applies to RQI. Such an amplification of unwanted small components does not play a role in DPA, because the vector to which the inversion is applied to is the same deflated vector  $\tilde{\mathbf{b}}$  in each step.

**5. Acceleration.** Subspace acceleration is very effective for solving linear systems (compare Richardson iteration and GMRES) and for eigenvalue computation (compare the power method and Arnoldi). Similar improvements are to be expected in the computation of dominant eigentriples [21]. The acceleration is an iterative process that, in each step, requires *expansion* of the search subspace (see §§5.1-5.2) and *extraction* of the approximate eigentriple from the search subspace (see §5.1). When the search subspace becomes too large (too high dimensional), some *restart* strategy is needed (cf. §5.4).

**5.1. Subspace acceleration.** Suppose  $\mathbf{X}$  is an  $N \times k$  matrix of which the columns span a search subspace for right eigenvectors. For computational convenience, suppose  $\mathbf{X}$  is orthonormal.

*Expansion.* To expand  $\mathbf{X}$ ,

$$(5.1) \quad \begin{aligned} &\text{Solve } (s_k \mathbf{E} - \mathbf{A})\mathbf{x} = \mathbf{b} \text{ for } \mathbf{x}. \\ &\text{Orthonormalize } \mathbf{x} \text{ against } \mathbf{X} \text{ to } \tilde{\mathbf{x}}. \\ &\text{Expand } \mathbf{X} : \mathbf{X} \leftarrow [\mathbf{X}, \tilde{\mathbf{x}}] \end{aligned}$$

Similarly, expand the matrix  $N \times k$  matrix  $\mathbf{Y}$  that spans the search subspace of left eigenvectors. The Gram-Schmidt process  $\mathbf{x} - \mathbf{X}\mathbf{X}^*\mathbf{x}$  orthogonalizes  $\mathbf{x}$  against  $\mathbf{X}$ . Then, normalization leads to  $\tilde{\mathbf{x}}$ , and we have *orthonormalized  $\mathbf{x}$  against  $\mathbf{X}$  to  $\tilde{\mathbf{x}}$* . In practice, a more stable variant (as repeated Gram-Schmidt) will be required.

*Extraction* [21]. The approximate eigentriple  $(\mathbf{x}_k, \mathbf{y}_k, s_{k+1})$  is computed from the ‘best’ eigentriple  $(v, w, s_{k+1})$  from the projected system  $(\mathbf{Y}^* \mathbf{A} \mathbf{X}, \mathbf{Y}^* \mathbf{E} \mathbf{X}, \mathbf{Y}^* \mathbf{b}, \mathbf{X}^* \mathbf{c})$ . Here, ‘best’ is the eigentriple with absolute largest residue  $(c^* v)(w^* b)$ , where  $c \equiv \mathbf{X}^* \mathbf{c}$  and  $b \equiv \mathbf{Y}^* \mathbf{b}$ . Then,  $\mathbf{x}_k = \mathbf{X}v$  and  $\mathbf{y}_k = \mathbf{Y}w$ .

Note that the expansion only requires the updated  $s_{k+1}$ . Nevertheless, having the approximate eigenvectors is useful too. They can be exploited in the stopping criterion. Obviously, they are needed when the approximations are sufficiently accurate.

The suggested extraction strategy, is tailored to finding dominant poles. The expansion strategy is the standard one for computing eigentriples.

**5.2. Orthogonalization or bi-E-orthogonalization.** As an alternative to the orthonormalization as suggested in (5.1), the expansion vectors  $\mathbf{x}$  and  $\mathbf{y}$  can also be ‘bi-E-orthogonalized’:  $\hat{\mathbf{x}} = (\mathbf{I} - \mathbf{X}\mathbf{Y}^* \mathbf{E})\mathbf{x}$  and  $\hat{\mathbf{y}}^* = \mathbf{y}^*(\mathbf{I} - \mathbf{E}\mathbf{X}\mathbf{Y}^*)$ . The resulting vectors can be scaled to  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$  such that  $\tilde{\mathbf{y}}^* \mathbf{E} \tilde{\mathbf{x}} = 1$ . If this strategy has been applied for obtaining all column vectors of  $\mathbf{X}$  and  $\mathbf{Y}$ , then  $\mathbf{Y}^* \mathbf{E} \mathbf{X} = \mathbf{I}$ .

This approach fits the fact that  $\mathbf{W}^* \mathbf{E} \mathbf{V} = \mathbf{I}$  (see §2.5), where  $\mathbf{V}$  and  $\mathbf{W}$  is the matrix with columns  $\mathbf{v}_j$  and  $\mathbf{w}_j$ , respectively. Another advantage is that the projected system (see (2.3)) reduces to a standard eigenvalue problem.

Unfortunately, the scaling may fail if  $\mathbf{y}^* \mathbf{E} \mathbf{x} = 0$ , or may be less stable if  $\mathbf{y}^* \mathbf{E} \mathbf{x}$  is small. Moreover, the inclusion of  $\mathbf{E}$  in the bi-orthogonalization makes the approach more costly.

**5.3. Why subspace acceleration.** The single vector iteration (DPA algorithm 3.1) offers only one choice for an approximate pole:  $s_{k+1} = \rho(\mathbf{x}, \mathbf{y})$ . Searching a subspace offers more possibilities for finding better approximations of the dominant pole: we select the eigenvalue with absolute largest projected residue  $(c^* v)(w^* b)$ . Note that, since  $\mathbf{x}$  and  $\mathbf{y}$  are in the span of the expanded  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively, the ‘DPA residue’  $(c^* \mathbf{x})(\mathbf{y}^* \mathbf{b}) / (\|\mathbf{x}\| \|\mathbf{y}\|)$  is in absolute value less than or equal to the selected ‘best’ residue. Therefore, we can state that acceleration leads to faster convergence and offers more control on the convergence: it enlarges the region of convergence of the dominant pole.

In addition, a better start is available after deflation, that is, after detection of one (dominant) eigentriple, we already have a space that will contain (good) approximations for the next (dominant) eigentriple. Note that it may be helpful to deflate the detected eigentriple not only

from  $\mathbf{b}$  and  $\mathbf{c}$  (as explained in §4) but also from  $\mathbf{X}$  and  $\mathbf{Y}$ . When combined with a restart strategy (see §5.4), these matrices have a few columns only, and deflation is not costly.

These three advantages (faster convergence, better control, better restart) compensate for the more expensive steps.

Block versions of DPA also employ subspaces (see [14]), but they miss the advantages mentioned above [21].

**5.4. Thick restart.** When the subspaces, that is, the  $\mathbf{X}$  and  $\mathbf{Y}$ , are expanded in each step, the memory requirements grow and the steps become computationally and increasingly more expensive. Therefore, a *restart* may become desirable, that is, a maximum acceptable dimension  $k_{\max}$  has to be selected and the search subspaces have to be reduced when they reach a dimension  $k$  larger than  $k_{\max}$ .

We suggest the following restart strategy.

Select a  $k_{\min}$ . Consider the projected  $k$ th order system

$$(5.2) \quad (\mathbf{Y}^* \mathbf{A} \mathbf{X}, \mathbf{Y}^* \mathbf{E} \mathbf{X}, b, c) \quad \text{with} \quad b \equiv \mathbf{Y}^* \mathbf{b}, \quad c \equiv \mathbf{X}^* \mathbf{c}.$$

If  $k > k_{\max}$ , then

- find the eigentriples  $(v_j, w_j, \mu_j)$  of (5.2)
- compute the associated residues  $(c^* v_j)(w_j^* b)$
- order the eigentriples such that the absolute value of the residues decrease
- continue with  $\mathbf{X} \leftarrow \mathbf{X}[v_1, \dots, v_{k_{\min}}]$ ,  $\mathbf{Y} \leftarrow \mathbf{Y}[w_1, \dots, w_{k_{\min}}]$ .

A *thick restart* [8, 29] with  $k_{\min} > 1$  is usually more effective than a complete restart where  $k_{\min} = 1$ . With  $k_{\min} > 1$ , we also keep an approximation in the search subspace for the eigentriple that is second in dominance. As explained in §5.3, this is helpful for starting the search to the next eigentriple. In eigenvalue computations, when an extraction strategy is based on Ritz-values closest to some target value, maintaining a search subspace of dimension larger than 1 is also helpful to diminish the effects of selecting a ‘wrong’ (ghost) Ritz-value (cf., e.g., [17, 8]). Since our extraction strategy here, is based angles between vectors (via the inner products  $\mathbf{c}^* \mathbf{X} v_j$  and  $w_k^* \mathbf{Y}^* \mathbf{b}$ ), there is less danger of misselection.

We have good experience with the values  $k_{\max} = 6$  and  $k_{\min} = 2$ .

**5.5. Subspace expansion.** As an alternative to the DPA expansion strategy as explained in §5.1, one can consider other expansion strategies, as RQI (cf. §3.3), but also JD (Jacobi-Davidson) [28].

Let  $(v, w, s_k)$  be the best eigentriple of the projected system (5.2), where, with  $\mathbf{v} \equiv \mathbf{X}v$  and  $\mathbf{w} \equiv \mathbf{Y}w$ , the vectors  $v$  and  $w$  are scaled such that  $\mathbf{w}^* \mathbf{E} \mathbf{v} = 1$ .

*Alternative expansion strategies.*

DPA: solve  $(s_k \mathbf{E} - \mathbf{A})\mathbf{x} = \mathbf{b}$  for  $\mathbf{x}$

RQI: solve  $(s_k \mathbf{E} - \mathbf{A})\mathbf{x} = \mathbf{E} \mathbf{v}$  for  $\mathbf{x}$

JD: solve  $(\mathbf{I} - \mathbf{E} \mathbf{v} \mathbf{w}^*)(s_k \mathbf{E} - \mathbf{A})(\mathbf{I} - \frac{\mathbf{v} \mathbf{v}^*}{\mathbf{v}^* \mathbf{v}})\mathbf{t} = \mathbf{r} \equiv (s_k \mathbf{E} - \mathbf{A})\mathbf{v}$  for  $\mathbf{t}$

Similar equations for expanding the search subspace for left eigenvectors have to be included.

JD solves for the (bi-E)orthogonal correction  $\mathbf{t}$  of the best available eigenvector approximations  $\mathbf{v}$ . The JD approach as compared to RQI has better stability properties: the system that is to be solved is better conditioned and it solves for a correction. Even when the correction is small, the JD expansion equation allows accurate computation, while in the RQI approach the small correction has to be computed as a difference of two relatively large vectors (two vectors that approximate the same eigenvector!).

We have the following surprising result (see also [27])

**THEOREM 5.1** ([22], Th.3.4.3). *With subspace acceleration, with exact solves of the expansion equations, no restart and appropriate initiation the three approaches produce the same approximate eigentriples in each step.*

The initiation of RQI and JD is related to initiation of DPA as indicated in Alg. 3.1 (see also §3.4):  $\mathbf{x}_0 = (s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}$  and  $\mathbf{y} = (s_0 \mathbf{E} - \mathbf{A})^{-*} \mathbf{c}$ . The initiation  $\mathbf{x}_0 = (s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{E} \mathbf{b}$  may seem to be more appropriate for RQI. However, in general, we often find in practise that  $\mathbf{E} \mathbf{b} = \mathbf{0}$  or  $\mathbf{E}^* \mathbf{c} = \mathbf{0}$ , and the more ‘appropriate’ initiation will fail then.

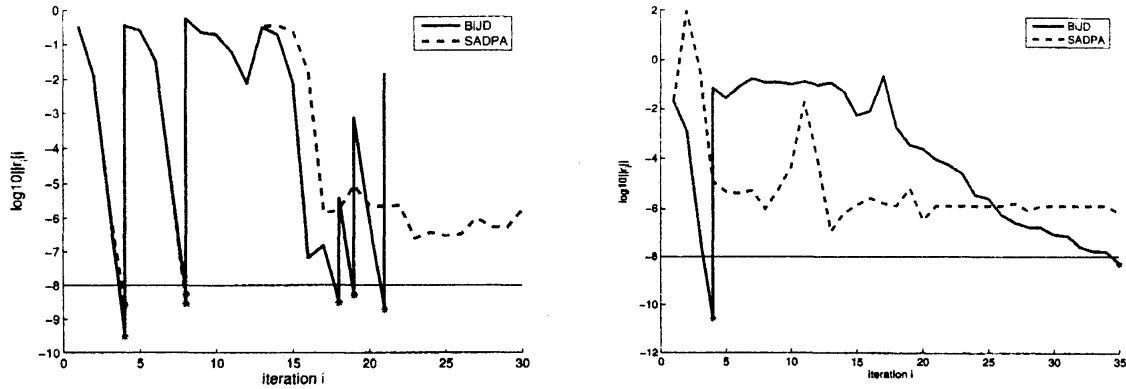


FIG. 5.1: The left panel shows the convergence history (the norm of the residual for the approximate right eigenvector) of subspace accelerated DPA (SADPA, the solid line) and JD (BiJD, the dashed line) for the PEEC model [6] when solving the linear systems exactly (with a new LU-decomposition for each system). The right panel shows the result for the Brazilian Interconnected Power System [21] when solving the linear system with only 10 steps of preconditioned GMRES. The same preconditioner (an LU-decomposition) is used in each step.

Recall that DPA allows efficient deflation, which is not the case for RQI and JD. This explains why subspace accelerated DPA is the method of choice when exact solves are affordable (that is, when dealing with systems of moderate order) [22, §3.3.3]. It should be noted, however, that DPA allows only limited accuracy. Because, if  $s_k$  does not change, then DPA does not lead to expansion. In this context, it is of interest to recall that the eigenvalue approximation is often much more accurate than the approximation of the eigenvector (since we approximate right and left eigenvectors at the same time, the error in the eigenvalue is proportional to the square of the error in the angle of the eigenvector):  $s_k$  can have full accuracy while the eigenvectors are still inaccurate.

If exact solves are not feasible, then JD is the method of choice [22, §3.6]. The system for expansion is better conditioned, and solves for a residual (small correction), it produces effective corrections even if  $s_k$  stagnates, preconditioners can be included. This is illustrated by the convergence histories in Fig. 5.1. The deflation, however, can become very expensive if many eigentriples are required (recall that the number of eigentriples is determined by the number of dominant eigentriples that are required for an accurate representation of the transfer function).

If, for instance, the expansion equation has been (inexactly) solved with a fixed number of Bi-CG steps, then subspace accelerated RQI and JD are still equivalent (in exact arithmetic, with appropriately matching initiation, no preconditioner, no restarts, see [27, 26]), but they are not equivalent to subspace accelerated DPA [22].

JD performs better in the left panel of Fig. 5.1 since it produces more accurate eigenvectors than subspace accelerated DPA. This is due to the fact that JD expands with corrections of the eigenvectors, rather than with approximate eigenvector eigenvectors as DPA does. The right panel shows the stabilizing effect of including projection in the expansion equations for JD.

**6. Modal approach or moment matching.** Moment matching (see §6.1 below) and modal approach are the two principal approaches for approximating the transfer function using models of reduced dimension. They have different effects. Moment matching tends to produce approximate transfer functions with small error in the ‘smooth’ part of the transfer function, while the error with the modal approach tends to be small in the high peaks and surroundings. Combining these two approaches can efficiently improve accuracy (see §6.3).

**6.1. Moment matching.** If  $s_0 \in \mathbb{C}$  is in the range of interesting values of  $s$  and  $s_0\mathbf{E} - \mathbf{A}$  is non-singular, then, for  $s$  close to  $s_0$ , we can form Neumann expansions: with  $\tilde{\mathbf{A}} \equiv (s_0\mathbf{E} - \mathbf{A})^{-1}\mathbf{E}$

and  $\tilde{\mathbf{b}} \equiv (s_0 \mathbf{E} - \mathbf{A})^{-1} \mathbf{b}$ , we have

$$H(s) = \mathbf{c}^* (s \mathbf{E} - \mathbf{A})^{-1} \mathbf{b} = \mathbf{c}^* (\mathbf{I} - (s_0 - s) \tilde{\mathbf{A}})^{-1} \tilde{\mathbf{b}} = \sum_{j=0}^{\infty} m_j (s_0 - s)^j,$$

where  $m_j \equiv \mathbf{c}^* \tilde{\mathbf{A}}^j \tilde{\mathbf{b}}$  are the so-called *shifted moments*. Projection onto the Krylov subspaces  $\mathcal{K}_k(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$  and  $\mathcal{K}_k(\tilde{\mathbf{A}}^*, \mathbf{c})$  leads to accurate approximations if the Neumann series converge quickly (if  $H$  is ‘smooth’ around  $s_0$ ).

For instance, if  $\mathbf{V}_k$  and  $\mathbf{W}_k$  are the ‘bi-Lanczos bases’ of these Krylov subspaces and  $T_k$  is the associated tridiagonal, then the first  $2k$  moments of the transfer function associated to the  $k$ th order system  $(s_0 T_k - I, T_k, \mathbf{W}_k^* \tilde{\mathbf{b}}, \mathbf{V}_k^* \mathbf{c})$  match the first  $2k$  moments of  $H$ . Here, for both functions, expansions are taken around  $s_0$ ,  $T_k$  is the  $k \times k$  upper block of the  $(k+1) \times k$  tridiagonal bi-Lanczos matrix  $\underline{T}_k$ :  $\tilde{\mathbf{A}} \mathbf{V}_k = \mathbf{V}_{k+1} \underline{T}_k$ . There are variants based on (two-sided), shift and invert Arnoldi that is, Arnoldi for the shifted and inverted operator  $\tilde{\mathbf{A}}$ .

Note that  $(s_0 T_k - I, T_k, \mathbf{W}_k^* \tilde{\mathbf{b}}, \mathbf{V}_k^* \mathbf{c})$  is a reduced model for  $(s_0 \tilde{\mathbf{A}} - \mathbf{I}, \tilde{\mathbf{A}}, \tilde{\mathbf{b}}, \mathbf{c})$  rather than for  $(\mathbf{A}, \mathbf{E}, \mathbf{b}, \mathbf{c})$ . Nevertheless, the matrices  $\mathbf{V}_k$  and  $\mathbf{W}_k$  can represent the interesting part of the spaces well, and the system  $(\mathbf{W}_k^* \mathbf{A} \mathbf{V}_k, \mathbf{W}_k^* \mathbf{E} \mathbf{V}_k, \mathbf{W}_k^* \mathbf{b}, \mathbf{V}_k^* \mathbf{c})$  may form a better and more natural reduced system. It often preserves properties as stability ( $\text{Re}(\lambda_i) < 0$  for all eigenvalues  $\lambda_i$ ) and passivity better.

The rational Krylov sequence (RKS, [25]) method (and two sided variants [11]) is a variant of shift and invert Arnoldi that allows to select different shifts  $s_0, \dots, s_m$  in each expansion step. The transfer function of the resulting reduced system tends to approximate  $H$  accurately in a neighborhood of each of the  $s_i$  (with accuracy depending on the number of iterates in which the same shift  $s_i$  is used). The neighborhood of accurate approximation tends to be larger if  $H$  is ‘smooth’ around  $s_i$ .

Fig. 6.1 (see [22]) shows approximation results for a configuration of the Brazilian Interconnected Power System [21]. The left panel shows the effect of working with a single shift  $s_0 = 0$  (and  $k = 85$ ) and with a second at  $s_1 = i$  ( $s_j = \sigma_j i$ ,  $k = 80$ ) using a rational Krylov Arnoldi (RKA) variant. The second shift is located at a peak. With one shift, we have an accurate approximation around the shift, but the good accuracy does not extend beyond the peak at  $s_1 = i$ . Including the peak as second shift extends the area of accurate approximation.

**6.2. Modal approach.** As observed in §2.5, the transfer function  $H$  tends to peak on values  $s$  along the imaginary axis that are equal to the imaginary part of dominant poles  $\lambda_j$ , with peaks being larger if the scaled residue  $R_j/\text{Re}(\lambda_j)$  is larger. The modal approach computes dominant poles and projects onto the spaces spanned by the associated eigenvectors. As a consequence, the peaks in the transfer function are extremely well represented by the transfer function of the reduced model (cf. §2.5).

**6.3. The effect of combining approaches.** In the left panel of Fig. 6.1, subspace accelerated DPA detected 36 dominant poles (20 poles in the upper half plane where required. Some poles show up in conjugate pairs, others are real). The *relative error* in the resulting approximate transfer function is displayed by the dotted line. A reduced model of order 18 was computed with RKA using shifts at 0, 2, 4, 6, 8, 10 times  $i$ , 6 moments for each of the shifts (that is, Krylov spaces of order 3). The solid line gives the resulting relative error. DPA detected the peaks (the dip in the dotted line at  $\omega = 1$ ), but RKA produces a better approximation away from the peaks (even with a smaller model). Note that neither of these two reduced models produce accurate approximations. Applying RKA to a model from which dominant triples (detected by DPA) have been removed, leads to a high accuracy in a wide range of  $\omega$  values: see the dashed line that corresponds to a reduced system of order  $k = 37$ . This system has been obtained by removing 19 poles (10 poles in upper half plane) with DPA and 18 RKA vectors. The dashed-dotted line represents the reduced model that uses 36 dominant eigentriples (as was the case for the dotted line) plus 18 RKA vectors (as for the solid line). The combination is much more accurate than will be anticipated from the ‘pure’ approaches. For more details on this discussion, see [22, §3.8].

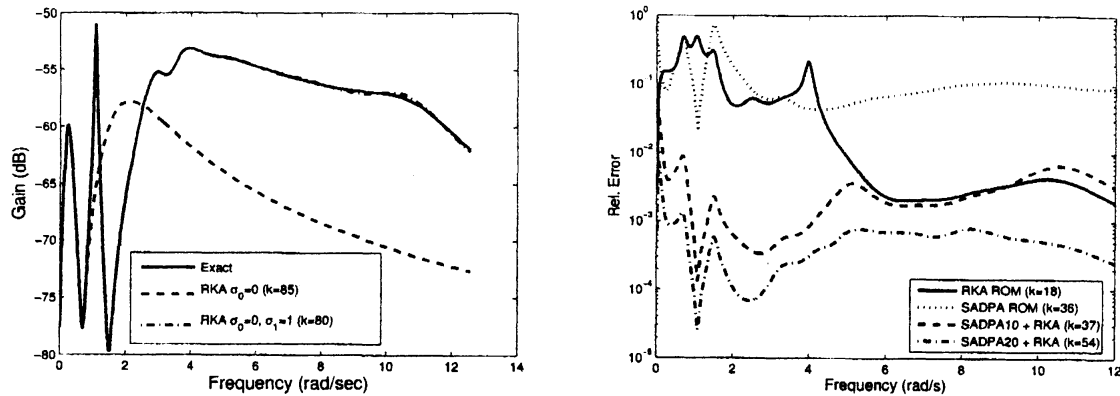


FIG. 6.1: The left panel displays the bode plot of the transfer function of the original system and of two reduced systems obtained with RKA, a rational Arnoldi variant. The right panel displays for the same system the relative error in approximations of the transfer function using subspace accelerated DPA, using RKA with six shifts, and using a combinations of these two methods (two combinations).

**7. Conclusions.** Models of reduced dimension can be formed by projecting the original system onto eigenvectors associated with dominant poles (the modal approach). The resulting approximate transfer functions tend to be accurate specifically around the peaky parts.

Although the Rayleigh quotient iteration (RQI) for computing eigentriples converges asymptotically faster than the dominant pole algorithm (DPA), the convergence regions of dominant poles tend to be much larger for DPA. This explains why DPA is more attractive than RQI for the modal approach for computing models of reduced dimension.

Inclusion of subspace acceleration in DPA and in RQI improves the efficiency of these methods. The accelerated versions of both DPA and RQI have the same convergence properties as Jacobi-Davidson (JD). In particular, the regions of convergence of all three methods coincide. Nevertheless, also when accelerated, DPA is more attractive than RQI (and JD): DPA allows efficient deflation of detected eigentriples. This makes the steps of DPA in the iterative search for subsequent eigentriples more efficient. All methods require the solution of linear systems. The observations so far refer to the case where the linear systems are solved exactly. For high-dimensional systems, this is not feasible. With inexact solves, JD has better stability properties and is the preferred method.

Moment matching techniques that project onto (rational) Krylov subspaces tend to produce approximate transfer function that are accurate on the smooth part. A combination with the modal approach for high accuracy of the peaky part of the transfer function can efficiently be performed and leads to high accuracy for a wide range of values of  $\omega$  (cf., (1.1)).

#### REFERENCES

- [1] A. C. Antoulas, *Approximation of large-scale dynamical systems*, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, USA, 2005.
- [2] A. C. Antoulas and D. C. Sorensen, *Approximation of large-scale dynamical systems: an overview*, Int. J. Appl. Math. Comput. Sci **11** (2001), no. 5, 1093–1121.
- [3] Z. Bai and Y. Su, *Dimension reduction of large-scale second-order dynamical systems via a second-order Arnoldi method*, SIAM J. Sci. Comput. **26** (2005), no. 5, 1692–1709.
- [4] C. Beattie and D.W. Fox, *Localization Criteria and Containment for Rayleigh Quotient Iteration*, SIAM J. Matrix Anal. Appl. **10** (1989), 80–93.
- [5] J. L. Casti, *Dynamical systems and their applications: Linear theory*, Academic Press, Philadelphia, PA, USA, 1997, Mathematics in Science and Engineering, Vol. 135.
- [6] Y. Chahlaoui and P. van Dooren, *A collection of benchmark examples for model reduction of linear time invariant dynamical systems*, SLICOT Working Note 2002-1, 2002.
- [7] E. Davidson, *A method for symplifying linear dynamic systems*, IEEE Trans. Aut. Control **11** (1966), 93–101.
- [8] Diederik R. Fokkema, Gerard L. G. Sleijpen, and Henk A. van der Vorst, *Jacobi-Davidson style QR and QZ algorithms for the reduction of matrix pencils*, SIAM J. Sci. Comput. **20** (1999), no. 1, 94–125 (electronic). MR 99e:65061

- [9] K. Glover, *All optimal Hankel-norm approximations of linear multivariable systems and their  $l^\infty$ -error bounds*, Int. J. Control **39** (1984), 1115–1193.
- [10] W. B. Gragg and A. Lindquist, *On the partial realization problem*, Lin. Alg. Appl. **50** (1983), 277–319.
- [11] E. J. Grimme, *Krylov projection methods for model reduction*, Ph.D. thesis, University of Illinois, Illinois, 1997.
- [12] A. M. A. Hamdan and A. H. Nayfeh, *Measures of modal controllability and observability for first- and second-order linear systems*, J. Guidance Control Dynam. **12** (1989), 421–428.
- [13] T. Kailath, *Linear systems*, Prentice-Hall, Engelwood Cliffs, NJ, 1980.
- [14] N. Martins, *The dominant poles eigensolver*, IEEE Trans. Power Syst (1997), no. 12, 245–254.
- [15] N. Martins, L.T.G. Lima, and H.J.C.P. Pinto, *Computing dominant poles of power system transfer functions*, Power Systems, IEEE Transactions on **11** (1996), no. 1, 162–170.
- [16] B. C. Moore, *Principal component analysis in linear systems: controllability, observability, and model reduction*, IEEE Trans. Aut. Control **26** (1981), no. 1, 17–32.
- [17] Ronald B. Morgan, *Computing interior eigenvalues of large matrices*, Linear Algebra Appl. **154/156** (1991), 289–309. MR 92e:65050
- [18] A. M. Ostrowski, *On the convergence of the Rayleigh quotient iteration for the computation of the characteristic roots and vectors. I, II*, Arch. Rational Mech. Anal. **1** (1958), 233–241 **2** (1958/1959), 423–428. MR 21 # 427
- [19] ———, *On the convergence of the Rayleigh quotient iteration for the computation of the characteristic roots and vectors. III. (Generalized Rayleigh quotient characteristic roots with linear elementary divisors)*, Arch. Rational Mech. Anal. **3** (1959), 325–340. MR 21 # 4541a
- [20] Beresford N. Parlett, *The Rayleigh quotient iteration and some generalizations for nonnormal matrices*, Math. Comp. **28** (1974), 679–693. MR 53:9615
- [21] J. Rommes and N. Martins, *Efficient computation of transfer function dominant poles using subspace acceleration*, IEEE Trans. Power Syst (2006), no. 21, 1218–1226.
- [22] Joost Rommes, *Methods for eigenvalue problems with applications in model order reduction*, Ph.D. thesis, Utrecht University, Utrecht, The Netherlands, June 2007.
- [23] Joost Rommes and Nelson Martins, *Efficient computation of multivariable transfer function dominant poles using subspace acceleration*, IEEE Trans. Power Syst (2006), no. 21, 1471–1483.
- [24] Joost Rommes and Gerard L. G. Sleijpen, *Convergence of the dominant pole algorithm and Rayleigh quotient iteration*, SIAM J. Matrix Anal. Appl. **30** (2008), no. 1, 346–363.
- [25] Axel Ruhe, *Rational Krylov: a practical algorithm for large sparse nonsymmetric matrix pencils*, SIAM J. Sci. Comput. **19** (1998), no. 5, 1535–1551 (electronic).
- [26] Valeria Simoncini and Lars Eldén, *Inexact Rayleigh quotient-type methods for eigenvalue computations*, BIT **42** (2002), no. 1, 159–182. MR 2003d:65033
- [27] Gerard L. G. Sleijpen and Michiel Hochstenbach, *Two-sided and alternating Jacobi–Davidson*, Linear Algebra Appl. **358** (2003), no. 1–3, 145–172.
- [28] Gerard L. G. Sleijpen and Henk A. van der Vorst, *A Jacobi–Davidson iteration method for linear eigenvalue problems*, SIAM J. Matrix Anal. Appl. **17** (1996), no. 2, 401–425. MR 96m:65042
- [29] Andreas Stathopoulos and Yousef Saad, *Restarting techniques for the (Jacobi-)Davidson symmetric eigenvalue methods*, Electron. Trans. Numer. Anal. **7** (1998), 163–181 (electronic), Large scale eigenvalue problems (Argonne, IL, 1997). MR 99j:65062